

# A Differentiated Services Implementation for High-Performance TCP Flows

Volker Sander<sup>\*†</sup> Ian Foster<sup>\*‡</sup> Alain Roy<sup>‡</sup> Linda Winkler<sup>\*</sup>

## Abstract

We examine the ability of the differentiated services (DS) quality of service (QoS) architecture to provide end-to-end QoS for high bandwidth TCP flows. Sliding window flow control makes TCP flows bursty, which can lead to packet losses; TCP's congestion control mechanisms mean that such losses have a large impact on achieved throughput. These problems are exacerbated for high-bandwidth, wide area flows where large windows are needed for performance. We report on experimental studies conducted on a DS testbed based on commercial routers; these studies allow us to identify DS configurations that can support multi-MB/s flows.

*Keywords: Differentiated Services, Quality of Service, Transmission Control Protocol*

---

<sup>\*</sup>Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL 60439, U.S.A.

<sup>†</sup>Central Institute for Applied Mathematics, Forschungszentrum Juelich GmbH, 52425 Juelich, Germany

<sup>‡</sup>Department of Computer Science, The University of Chicago, Chicago, IL 60637, U.S.A.

# 1 Introduction

Network Quality of Service (QoS) can be addressed by a variety of mechanisms. However the actual deployment of those mechanisms in the current structure of the Internet is rare. Network layer mechanisms such as ATM assume a homogeneous infrastructure between two end-systems, which for economic reasons may not exist. The Integrated Service (IS) [3] Architecture does not depend on homogeneous infrastructure but often requires that the end-systems have a specialized kernel installed. Furthermore it requires a per flow handling on every networking device which results in scalability problems. Recently the IETF has defined a new QoS framework that addresses these concerns: The Differentiated Services (DS) [2] Architecture.

The DS architecture defines an architecture for implementing scalable service differentiation in the existing Internet. Edge routers are configured to mark packets for a desired class. Interior routers treat packets based on their aggregates (classes) which eliminates the scalability problem introduced by the IS model. The DS model defines different ways [12, 10] in which a marked packet should be treated by an interior router, based on the associated class. This specification is called per-hop behavior (PHB). In summary, an implementation of the DS architecture has to provide edge routers which are capable in marking packets per flow and interior routers which implements the related PHBs.

Several recent studies [5, 21, 22] show that it is difficult to guarantee requested throughput for TCP flows, due to the flow and congestion control [19, 1] mechanisms used by this protocol, which result, for example, in bursty traffic. As a result of these studies, researchers have proposed both changes to the DS architecture and modifications to the TCP protocol (e.g., Two-window TCP [5]). However, no study has analyzed high-speed flows (10 MB/s or more) when using commodity router hardware which provides buffer capabilities of up to several Megabytes. Distance visualization applications are just one sample of applications encountered in science and engineering which involve data transfers and media streaming at hundreds of megabytes per second (MB/s). Deploying QoS to end-systems requires a detailed analysis of actual mechanisms provided by commodity networking hardware, and their impacts on TCP performance.

Using Cisco's implementation of DS mechanisms, we performed a careful evaluation of high-bandwidth TCP performance and draw conclusions about how to configure DS to provide a premium service in this situation. This work was performed in the context of the Globus Architecture for Reservation and Allocation (GARA)[7], which builds on these mechanisms to deliver per-flow, advance reservation, end-to-end Quality of Service.

## 2 The Differentiated Services Architecture

### 2.1 Overview

The Differentiated Services (DS) architecture [2] is a reaction to earlier per-flow based network quality of service (QoS) architectures such as RSVP. In contrast to these per-flow, end-to-end QoS architectures, DS focuses on defining the behavior of aggregates. Packets are identified by simple markings which indicate which aggregate behavior they should be given. In the

core of the network, routers need not determine which flow a packet is part of, only which aggregate behavior should be used.

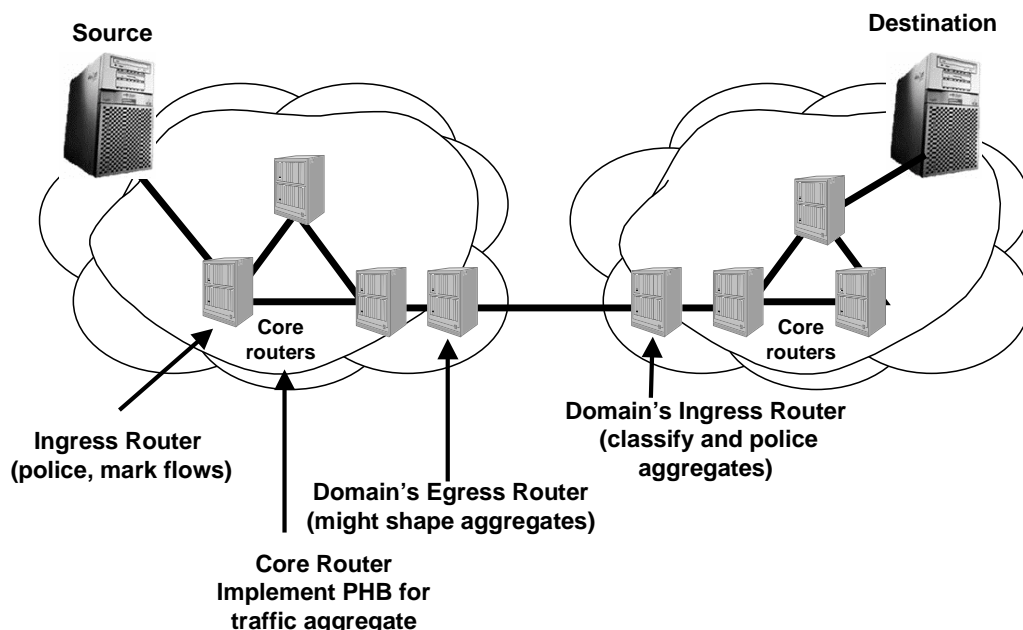


Figure 1: A network flow going through two DS domains. The ingress router is responsible for policing and marking of packets. The interior router handle the traffic based on the aggregate's PHB. The domain's egress might be enforced to shape the traffic of the aggregates. The domain's ingress router does policing based on aggregates.

DS pushes complexity to the edges of the network: packets are marked either by applications or by edge routers. If edge routers mark packets, they may choose to do so on a per-flow basis, or on any other criteria. Packets may be marked only when they are “within profile”—that is, when the sender is sending within predetermined limits such as bandwidth or time of day. In contrast to this complexity in the edge routers, routers in the core of the network provide service based only on these markings. A particular marking on a packet indicates a *Per-Hop Behavior* (PHB). Currently, the IETF's Differentiated Services Working Group has specified multiple PHBs [12, 10]. Within this paper, we consider the Expedited Forwarding (EF) PHB. EF is for high-priority packets, and routers are configured to allow up to a certain portion of the bandwidth on a link to be given to packets marked with EF,

and extra packets are dropped.

Figure 1 illustrates the DS architecture in a two domain environment.

## 2.2 Cisco's Implementation of Differentiated Services

We use a DS implementation based on Cisco 7500 series routers, which support the EF PHB, via two mechanisms:

- Committed Access Rate (CAR) is used on the ingress ports of edge routers to mark and police the flows for which premium bandwidth is required. It is also used on the ingress router of a domain to police the premium aggregate. CAR currently does not support specific DS packet markings (as defined in [16]) but instead uses the IP Precedence field to mark packets belonging to EF flows.
- Weighted Fair Queuing (WFQ) is used on the egress port of edge routers and in interior routers. WFQ ensures that in periods of congestion—i.e., when packets get queued in the router because the output link does not provide the capacity for delivering them immediately—each IP precedence class receives at least the fraction of the output bandwidth defined by the weight defined for that class. Cisco's VIP2-50 implementation allows premium classes to be assigned up to 99% of the available link bandwidth. It is important to note that WFQ only has an active effect when there is congestion. When the interface is not congested, queues can use any available bandwidth, regardless of their weight.

## 3 TCP Overview

Establishing end-to-end QoS for high-performance TCP flows in a DS environment is a challenge. This section gives an overview of the dynamics of TCP flows, and discusses the nature of associated challenges.

### 3.1 Flow Control

TCP implements flow control mechanisms to prevent a sender from transmitting data faster than its receiver can handle. These mechanisms constrain the maximum amount of outstanding data (i.e., data sent but not acknowledged) to be less than or equal to a specified window size. The sliding window protocol [20, 4] used to implement flow control uses two windows:

- The receiver advertises the offered or advertised window which corresponds to the amount of free space in the receiving buffer, relative to the acknowledged data. The receiver guarantees that it is able to buffer the data indicated by the window. The maximum offered window size can be set by the application.
- Based on the offered window size and on other information (e.g. the available local buffer space, the first unacknowledged segment, and the outstanding segments), the sender computes the usable (i.e., the actual) window size.

## 3.2 Congestion Control

Wide area tests have shown that if a sender immediately injects multiple segments into the network, up to the window size advertised by the receiver, the application throughput, often called goodput, can be reduced drastically [11] in the presence of competing traffic (congestion).

Hence, TCP has been modified to avoid the injection of a full offered window at the beginning of a session. The basic idea is to start with a small usable window size and then increase this window size based on the successful receipt of segments. Another TCP sender window is introduced, the congestion window (cwnd), which can reduce the usable window size—i.e., it is the minimum of the previous value and the cwnd. The operation of the sender is split into two phases:

- Within the slow start phase the cwnd is increased with every acknowledged segment, until it either reaches a system wide maximum threshold, or the transmitter detects that packets were dropped <sup>1</sup>. Starting with a cwnd-size of 1 or 2, the slow start phase doubles the cwnd-size with every round-trip of a full window.
- Within the congestion control phase the cwnd is only increased by one with every full window round-trip. Again this phase is limited by a system-wide threshold, and by the detection of dropped packets.

Figure 2 gives a more formal overview of the window evolution.

To emphasize the effect of these mechanisms, it is important to note that the goodput of a TCP application is limited to the actual window size, divided by round-trip time (bandwidth\*delay product [20]). The impact of a single drop on the short term TCP goodput is enormous. The macroscopic behavior of these mechanisms are described in [14, 13]. The general impact of the window size on the goodput is evaluated in [17, 15]. Statistical analysis for the impact on a DS implementation is discussed in [18, 22].

## 4 Experimental Studies

We report on experiments designed to examine a DS implementation based on commodity networking products.

### 4.1 Experimental Configuration

Our experimental configuration, illustrated in Figure 3, comprises a laboratory testbed at Argonne National Laboratory (the Globus Advance Reservation Network Testbed: GARNET) connected to a number of remote sites, including Lawrence Berkeley National Laboratory (LBNL). Connectivity to LBNL is provided by the Energy Sciences Network (ESnet) DS testbed. GARNET allows controlled experimentation with basic DS mechanisms; the wide

---

<sup>1</sup>Dropped packets are either determined by a time out or by receiving a specific amount of duplicate acknowledgments ([19, 1])

Initialize  $W = 1$  (or 2) and  $T$  with `max_cwnd`

- (1) After every non-repeated ACK:  
if  $W < T$ , set  $W = W + 1$ ; Slow Start Phase  
else set  $W = W + 1/[W]$ . Congestion Avoidance Phase
- (2) When the number of repeated ACKs exceeds a threshold  
(fast retransmission/fast recovery) retransmit "next expected" packet;  
set  $T = W/2$ ;  
set  $W = T$ ; (i.e. halve the window)  
resume congestion avoidance using the new window
- (3) Upon time expire, the algorithm goes into slow start:  
set  $T = W/2$ ;  
set  $W = 1$ .

Figure 2: Description of TCP-Reno's window size evolution (see [13] for details). Here,  $W$  represents the congestion window size and  $T$  is a threshold used to switch between the different phases, often referred as *ssthresh*.

area extensions allow for more realistic operation, albeit with a small number of sites. Notice that end-system resources are located in different domains; hence, we must deal with distributed authentication and authorization.

Cisco Systems 7507 routers are used for all experiments. Within GARNET, these are connected by OC3 ATM connections; across wide area links, they are connected by VCs of varying capacity. We are restricted to these relatively slow speeds because the 7507 cards do not implement CAR and WFQ at speeds faster than OC3. End system computers are connected to routers by either switched Fast Ethernet or OC3 connections.

## 4.2 Evaluation Tools

The analysis required appropriate measurement methods. It was necessary to use a tool which was capable of producing a network flow (TCP and/or UDP) with specified characteristics, such as bandwidth or burstiness. Furthermore there was the need to find a tool for creating a well-defined amount of competing traffic, to create some congestion on the internal network.

- *TCP stream generator*: we use a TCP traffic generator capable of generating a flow with a predetermined rate. This generator operates by adapting the frequency of `write()` calls to achieve the desired rate. Note that this only adapts the frequency with which the transmitter fills the socket buffer, and not the actual transmission rate. The host's TCP stack is still responsible for transmitting those packets onto the network media, by using TCP's flow and congestion control mechanisms. However, this behavior reflects exactly the challenge of QoS-aware TCP communication.

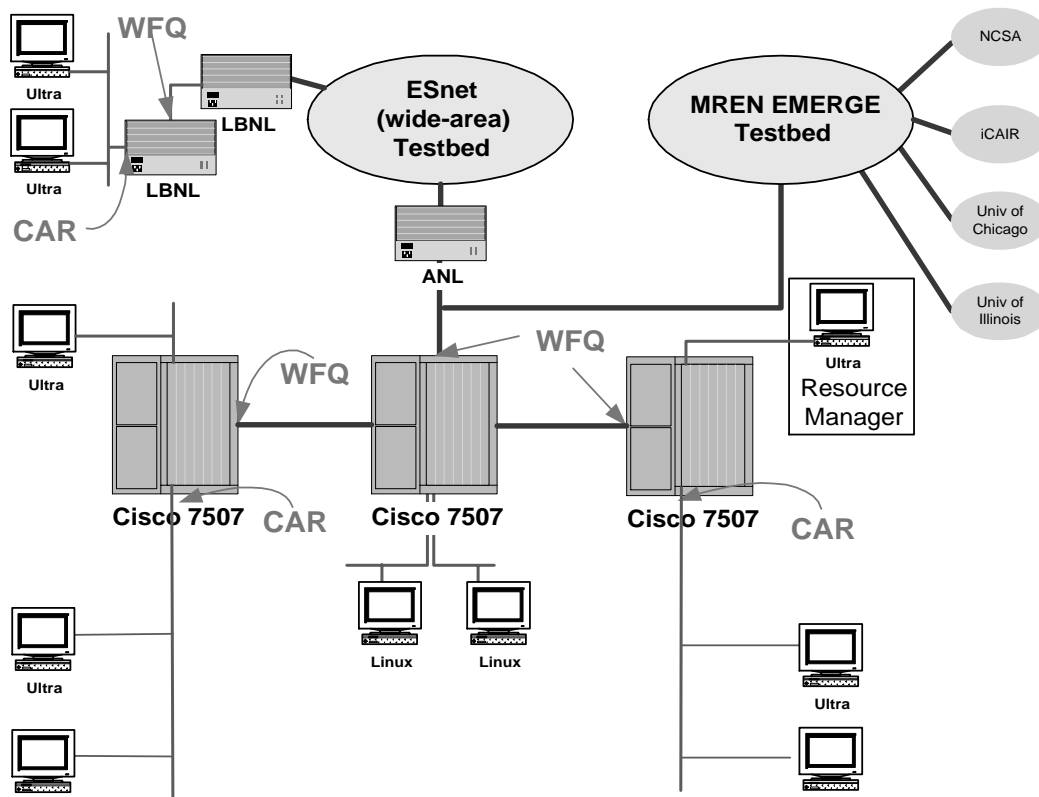


Figure 3: The experimental configuration used in this work, showing our local GARNET testbed and its extensions to remote sites connected via ESnet and MREN.

- *UDP traffic generator*: The basic idea of the traffic generator used (Andy Adamson, University of Michigan) is to divide the flow into intervals of one second, in which UDP datagrams are transmitted at fixed frequency until a given per second rate was achieved. Having achieved the desired amount of data per second, the transmission stops until the next interval starts. If the rate was not achieved, the size of the datagram used for the next second is increased by two bytes. If the sender has already submitted this amount of data for a second, it stops transmission until the next interval starts.

### 4.3 Evaluation of CAR

Our first experiments are designed to evaluate the impact of the policing performed by Committed Access Rate (CAR) on TCP flows. To this end, we ran flows over an otherwise uncongested link, with CAR applied at the ingress point.

### 4.3.1 Short-term Transfers

As stated above, TCP reacts when dropped packets are detected. If the transmitter recognizes a lost packet, it either falls into congestion control phase, or slow start phase, depending on the number of contiguous packets lost. This behavior has a dramatic impact on TCP performance. The first experiment (Figure 4) was designed to illustrate this impact, by demonstrating the behavior of a relatively small TCP transfer (16.7 Mbytes), with different transmission rates in combination with CAR.

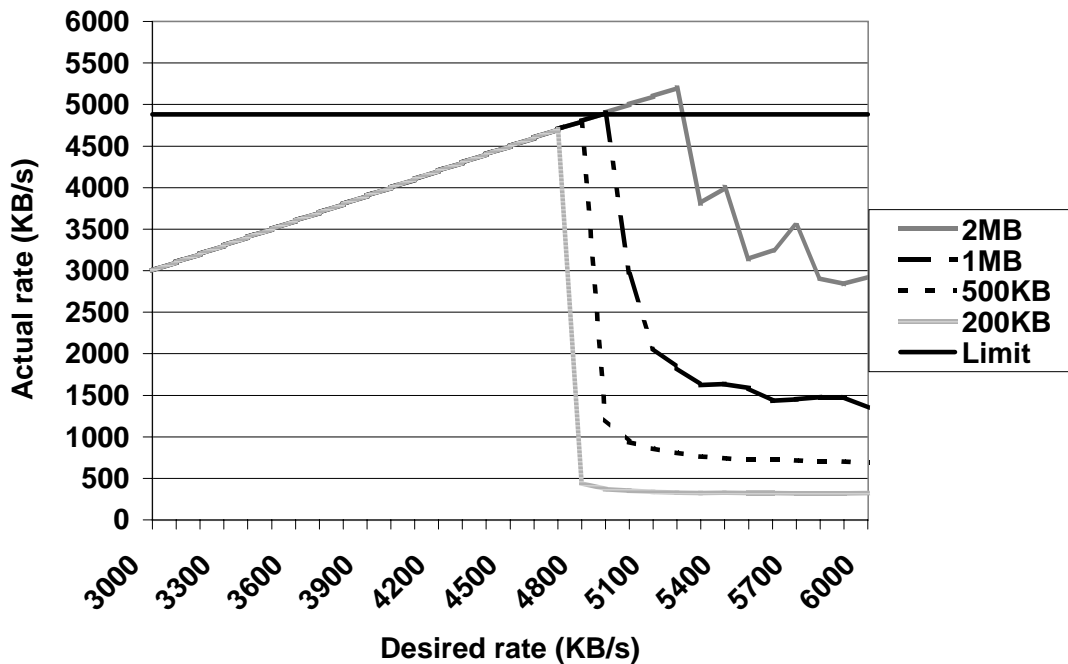


Figure 4: Received throughput of different short term TCP streams with an active CAR configuration on the ingress site. The token bucket depth was varied between 200 KB and 2 MB. The average rate limit used for configuring CAR was constant.

The results displayed in Figure 4 were gathered over several sessions by transmitting a TCP stream at a pre-specified rate and altering the CAR configuration. The impact of the token bucket depth<sup>2</sup> is recognizable. Increasing the normal burst size results in a higher

<sup>2</sup>Cisco has implemented the rate policies with a token bucket mechanism. The depth specifies the maximum amount of tokens a token bucket can store.



short-term bandwidth, until CAR's exceeding policy comes into play. Note that the increase of the bandwidth from 4700 KB/s to 5200 KB/s corresponds to the transmission time and the difference in the bucket depth. Having additional tokens for 1.8 Mbytes allows a stream to exceed its limit by 500 KB/s for nearly 4 seconds. Overall the transmitter is transferring 16.7 Mbytes which also can be done in less than 4 seconds. For an exact calculation it should be mentioned that the transmitter was using the path MTU discovery. For that reason the used segment size was 1460 Bytes.

However the above trial demonstrates effectively the behavior of TCP's slow start feature. As soon as the token bucket is empty, the router starts dropping packets and will often drop consecutive packets. The TCP transmitter is recognizing that consecutive packets are lost and reacts to this by shrinking the congestion window to two, which has an enormous short-term impact on the actual throughput.

From this, one can conclude that exceeding the actual rate limit causes TCP to decrease performance dramatically. Transmitting packets faster than the QoS rate limit has a negative impact on overall performance achieved. Note that TCP might submit a whole socket buffer in one burst as permitted by the offered receiver window and the congestion window. The actual CAR configuration must be able to handle those bursts without dropping packets. Currently the normal burst size is limited to 2 Mbytes, this has the side effect of shrinking the maximum supported window size to 2 Mbytes. Cisco has mentioned that future versions of CAR might be able to handle deeper token buckets.

#### 4.3.2 Long-term Transfers

Besides the short-term behavior of TCP streams, it is important to analyze a stream with a longer duration. For that reason several long term tcp sessions were monitored. Figure 5 represents an example for the behavior of a TCP session exceeding its rate limit.

In this case the normally constant throughput starts oscillating. As soon as CAR starts dropping packets, the transmitter reacts with TCP's slow start feature. This reduces the transmission rate drastically. As the router's token bucket begins filled again, the transmitter is able to exceed its limits again (if enforced by the application), which will resume the oscillating behavior.

The conclusion is that exceeding the actual rate limit causes long-term TCP traffic to oscillate. This is mainly because of the slow start/congestion avoidance feature and the bucket depth. The variation increases the greater the bucket depth.

### 4.4 Evaluation of WFQ

As stated before, Weighted Fair Queuing (WFQ) is conceptionally well suited for supporting the IETF's EF PHB. This section describes experiments done to evaluate the functionality of WFQ. Because we are focusing on end-to-end QoS, we were more interested in analyzing whether a specific TCP goodput of a high-end application could be guaranteed under congestion than analyzing the jitter of a UDP stream.

The test suite for creating well-defined congestion was the following:

There was a TCP stream trying to send a unidirectional stream of 10 MB/s. The edge routers were configured to mark packets of that stream and limit its rate to 11 MB/s. To

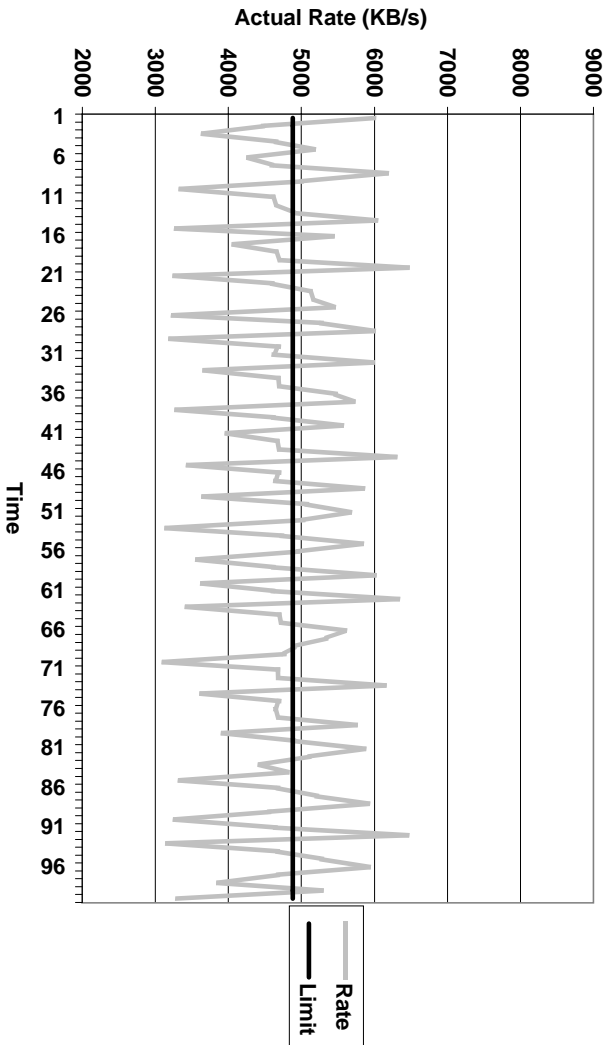


Figure 5: Received throughput of a TCP stream with an underestimated active *CAR* configuration on the ingress site. The attempted transmission rate was 6 MB/s. *CAR* was configured with a bucket depth of 2 MB.

emulate the behavior of a wide area TCP application we used a window size of 1 MB. Note that we were also required to increase the maximum congestion window size of our Solaris 7 end-system to 1 MB, because the default limitation is 262000 bytes.

After the TCP stream was running for a while, a UDP stream was started between the other SUN workstations, across the same link. Every 0.1 msec, packets of 500 bytes were submitted, which resulted in an application bandwidth of 5 MB/s.

To create congestion, we shrank the PVC between the interior router and the egress router to 100 Mbps. Note that this bandwidth includes all overheads of an ATM link, which means that the actual supported IP bandwidth is much lower, depending on the actual packet size. We did several different experiments:

#### 4.4.1 Impact of Congestion on a TCP Flow without QoS

The first experiment was done to demonstrate the behavior of commodity routers without WFQ. We examined this non-DS case in order to compare it with the DS case when there are no DS flows, because we want to be sure that the DS mechanisms do not impact the best-effort traffic.

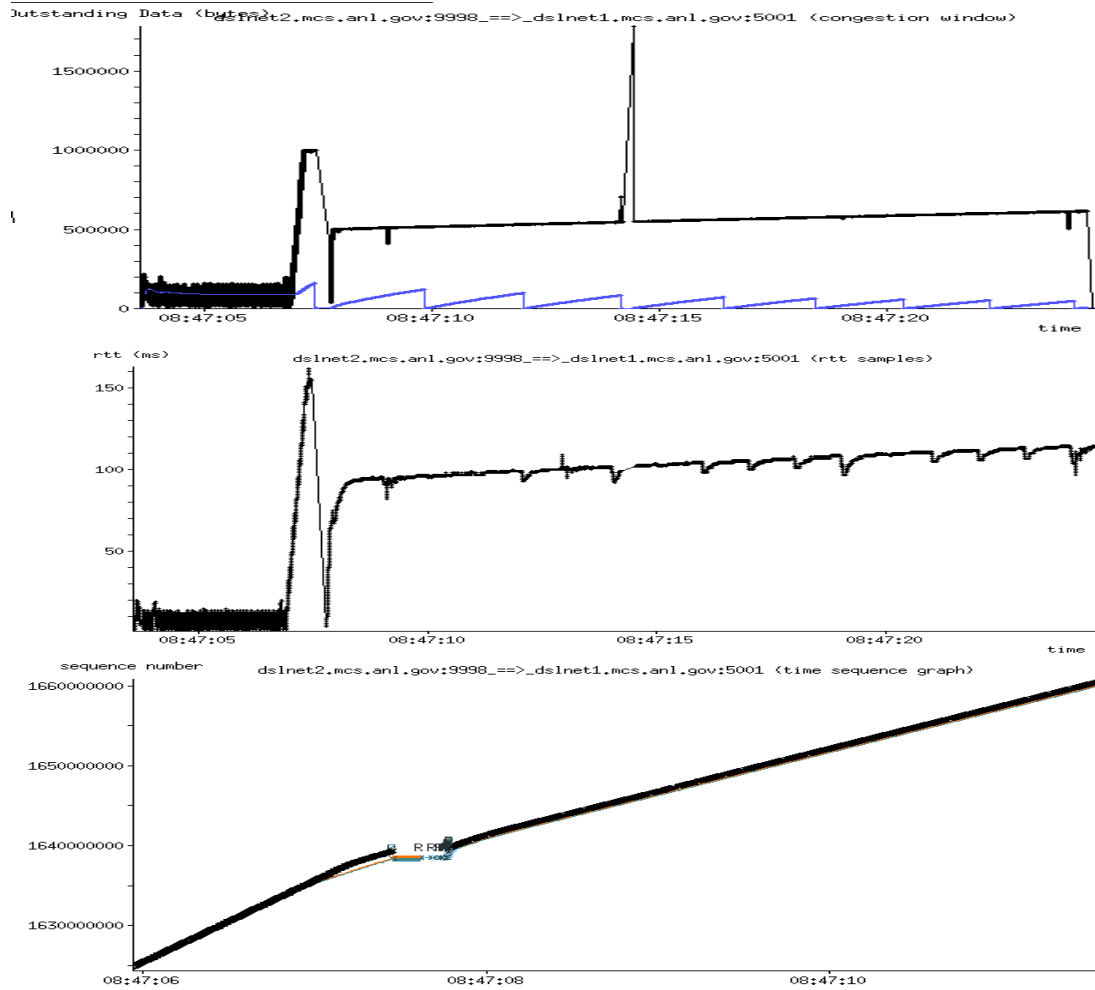


Figure 6: Performance achieved for a TCP flow under congestion on GARNET. We demonstrate the behavior of our experiment without any WFQ configuration. See text for details.

Figure 6 shows three graphs:

- The amount of outstanding data is displayed as upper line. It indicates the amount of packets which are not yet acknowledged. At 08:47:07 the congestion started. The amount of outstanding data immediately increased to its limit of 1 MB, thereafter it decreased to 500,000 – 600,000 Bytes. The peak at 08:47:15 was caused by a miscalcu-

lation due to CPU scheduling. It has no corresponding entry in the round-trip graph described below.

- The graph of the round-trip time (rtt) needs to be read carefully. TCP implementations do not acknowledge every received segment. The delayed acknowledgment mechanism delays this until the acknowledgement can be piggy-backed, or a timer expires (50 ms in the testbed). For that reason, the minimum value of this graph represents the actual round-trip time, but the maximum value indicates the value which is of interest for calculating the bandwidth-delay product for optimizing the window size. The graph itself correlates to the outstanding data. Without congestion the actual round-trip time is as low as expected. As congestion starts, the router starts queueing packets. This increases the round-trip time enormously. Under congestion the round-trip time is still extremely high. As expected, TCP does almost no piggy-backing of acknowledgements, because the actual transmission time is larger than the timeout of 50ms.
- The third graph shows the transmitted sequence numbers over time. Note that the time-scale of this graph is zoomed in, to clarify the actual behavior. The gradient of the graph shown indicates the actual goodput. The important thing to note is that at the time the congestion occurred, the graph lists a lot of “R”s. This means that the transmitter has retransmitted the same sequence number. Though we used the Selective Acknowledgment (SACK) feature, the number of retransmits listed in the graph indicates that the TCP flow was going into the congestion control phase and eventually into the slow start phase. The consequence of this is that the transmitter was backing off its transmission speed drastically. One can see that TCP went to slow start mode because the amount of outstanding data was reduced to a very small value (two segments). Now there was no longer congestion (the UDP traffic was only submitting 5 MB/s). As long as TCP was in slow start mode, it increased its window size exponentially. Having reached half of the former window size of 1 MB, it went to congestion control phase. Now the window size was only increased linearly. As a consequence the throughput was increasing permanently, but only by a small amount. This created more and more congestion on the network which lead to a higher round-trip time which slowed down the increase of the window size permanently.

The picture demonstrates that commodity networking hardware does have a significant amount of buffer capacity. It is not a problem for the router to buffer more than 500 KB of data. For that reason it is quite important to evaluate the impact of a WFQ configuration emulating the IETF’s EF PHB.

#### **4.4.2 Impact of Congestion on a non-premium TCP Flow with Standard WFQ Buffer Limits**

The following graph show the same experiment, without any premium traffic. But this time WFQ was configured to provide 98% of the output bandwidth to the premium class during congestion. The goal was to demonstrate the impact of a standard WFQ configuration on best effort traffic. Our hope was that enabling WFQ would not have any impact on best-effort traffic.

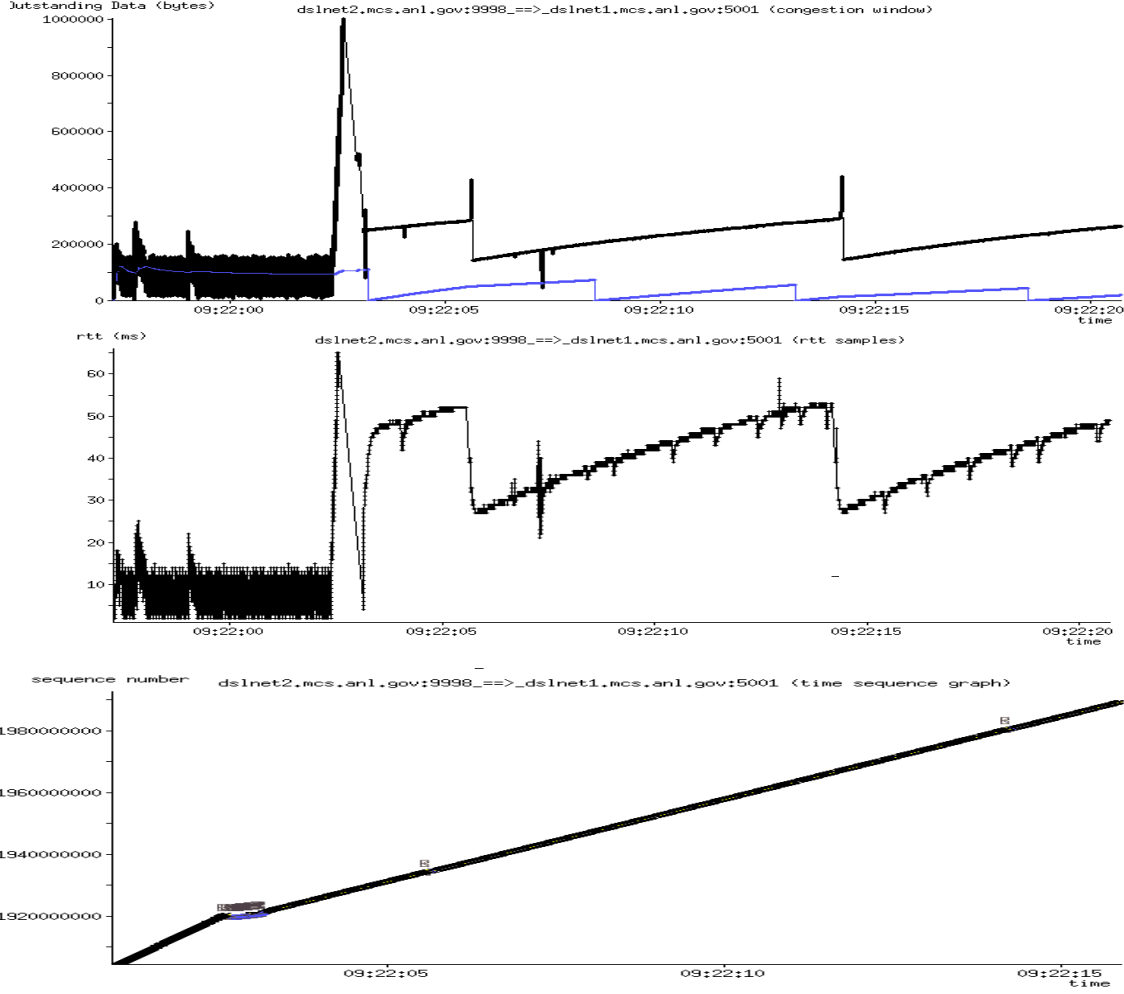


Figure 7: Performance achieved for a TCP flow under congestion on GARNET. We demonstrate the behavior of our experiment without any WFQ configuration. See text for details.

Figure 7 shows three graphs:

- The amount of outstanding data is different with WFQ. Whereas the configuration without WFQ was able to keep a 1 MB window outstanding over several hundreds of msec, this setup was not able to handle this burst. Though the number of outstanding bytes quickly reached the 1 MB limit, the actual amount of data buffered correctly was less than 550,000 bytes. This is indicated by the amount of outstanding data in the following congestion control phase, which is half of the former congestion window size. Even an amount of 550,000 bytes was not handled properly over the time. Because TCP was now increasing the actual window size only linearly—before it went to slow start mode again it was increasing it exponentially—it gives a better overview about the amount of data buffered correctly. In this case, it falls into congestion control phase again at 300,000 Bytes, which reflects the amount of buffer space available under this

condition.

- For the round-trip time results are similar to the amount of outstanding data. Because of the smaller amount of buffer space, the actual round-trip time is significantly less. The maximum round-trip time handled correctly is approximately 52 ms. Assuming that the 100 Mbps PVC goodput capacity of 86 Mbps (after SONET, ATM and protocol overhead), and assuming that the amount of time the acknowledgment packet needs is 2ms, the overall amount of buffer space is around 530,000 Bytes. Nearly two-third of this buffer space is available to the TCP flow which results in 350,000 bytes. The remaining difference between this rough calculation and the observed amount is caused by the fact that the UDP packets contain less packet header and by the fact that the WFQ implementation handles its memory management in chunks of 512 Bytes. This causes an internal fragmentation of the available queue space.
- The transmitted sequence number graph indicates that the activated WFQ configuration drops more packets than a router without WFQ configured. Again it needs to be noted that this graph was zoomed.

The reason for this behavior is that a WFQ configuration divides nearly the whole available amount of queue space to the aggregate specific queues based on their weight. For that reason a best effort aggregate with a guarantee of 2% of bandwidth under congestion receives only a small amount of queue space. Following the idea of a premium service, it does not make sense to provide a significant amount of queue space to the premium flow, because it would result in a potentially higher latency. Consequently an implementation of the EF PHB should change the default buffering drastically by providing a significant amount of queue space to the best effort queue.

A proper configuration—though in absence of an available traffic shaping mechanism (distributed traffic shaping was not yet supported) it was still providing a an amount of queue space for the premium class—could reproduce the same behavior than it was observed with the non-WFQ configuration even with assigning 99% to the premium queue <sup>3</sup>.

#### 4.4.3 Impact of Congestion on a Premium TCP Flow with Correct Configuration

Finally we want to demonstrate that the proposed implementation of the EF PHB really provides a premium channel. WFQ guaranteed 99% of the available bandwidth to the premium class. For test purposes, the policy on the ingress site was allowing 90% of the available rate. Note that in a real environment the absolute amount of allowed premium bandwidth is less.

Figure 8 shows three graphs:

- The amount of outstanding data is only increasing slightly when congestion started at 11:16:08. Actually the increase of outstanding data is less than 10%. Note that the

---

<sup>3</sup>Due to unavailability of equipment we could not finish this result, but a final version of a paper will contain this result

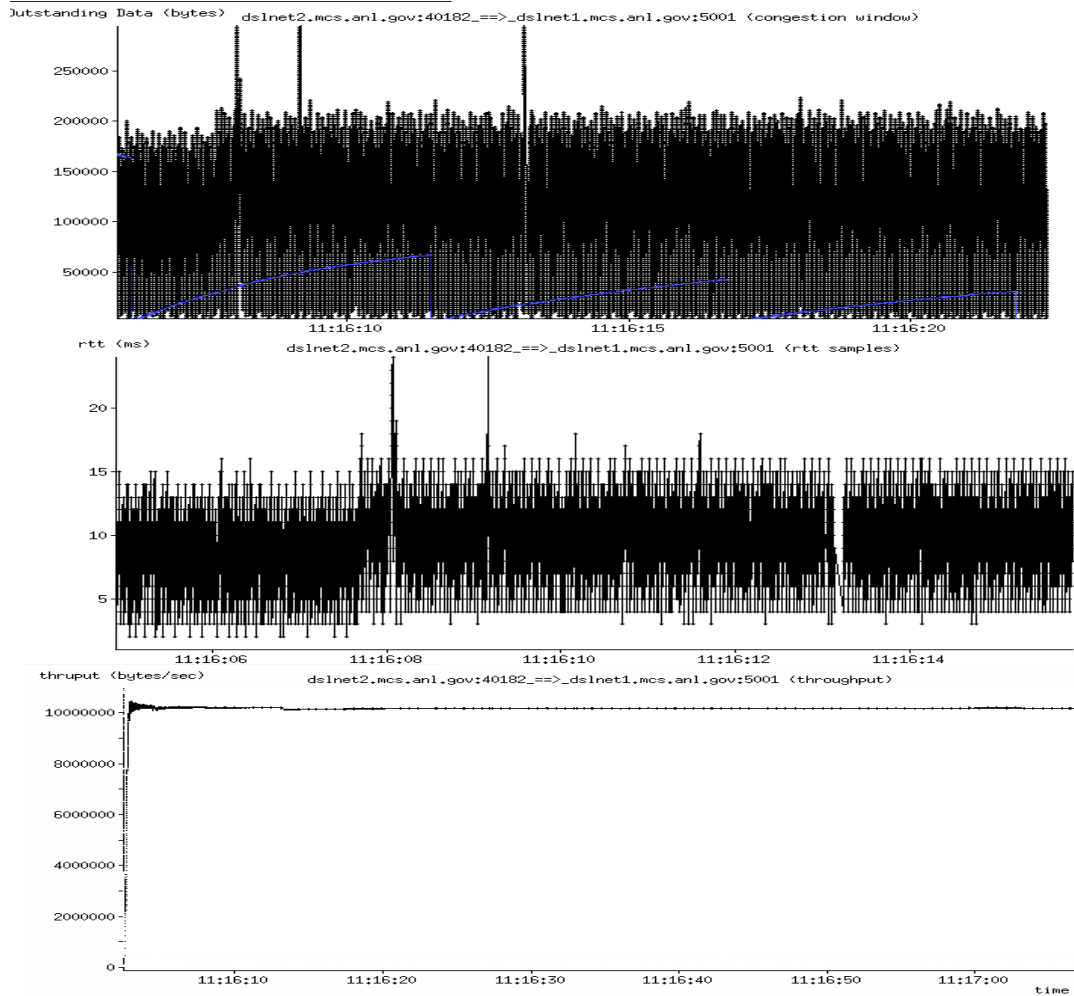


Figure 8: Performance achieved for a TCP flow under congestion on GARNET. We demonstrate the behavior of a premium high bandwidth TCP stream. See text for details.

scale of this graph differs from the previously shown outstanding data graphs, because this configuration does not introduce a buffering of segments inside the router.

- The round-trip time of the premium traffic is nearly not affected by the existence of congestion. Note that this graph is slightly zoomed. We see only a minor increase from 3ms to 4ms which is also result of the granularity of graphing tool, which is only able to display round-trip times in units of ms.
- The third graph displays the actual average throughput over the time. As we can see it is not influenced by the congestion.

## 5 Proposed Differentiated Services Configuration

Based on the experiments we propose the following DS configuration:

CAR and WFQ should be used for QoS enforcement, as described above. Flow specifications supplied to CAR should use a bandwidth computed from the user-specified required bandwidth, taking into account packet headers (note that this requires packet size information), and packets that exceed the rate-limit should be dropped. The burst size parameter should be set to the bandwidth (in bytes/second) multiplied by the assumed maximum round trip time.

To implement the EF PHB without traffic shaping, we recommend assigning 99% of the available bandwidth of the slowest link ( $BW$ ) to the premium class of WFQ. This is to avoid queuing of premium traffic as much as possible. Reducing this high amount of guaranteed bandwidth would not prevent any link from submitting 99% of premium traffic in case of no congestion, because WFQ scheduling only affects the behavior under congestion. This may seem like a large percentage to give to premium traffic, but CAR will ensure that premium traffic is properly limited. We will show that this configuration minimizes the queuing time for packets.

All other interior interfaces should exactly guarantee 99% of the rate  $BW$ . In doing this, we can now calculate the maximum queuing time introduced for the premium traffic. Let  $x * BW$  be the amount of premium bandwidth allowed by the policy (e.g, the average rate limit of CAR), where  $0 < x < 1$ . Furthermore let  $r_{tt}$  be the round-trip time with no competitive traffic and  $q_{t}$  the maximum queueing time for the premium traffic, introduced by bursts. Now we can calculate the maximum estimated round-trip time used by CAR to calculate the token bucket depth with  $r_{tt} + q_{t}$ . Assuming this round-trip time CAR actually limits the maximum size of a burst to  $x * BW * (r_{tt} + q_{t})$ . Because we can guarantee a premium bandwidth of  $0.99 * BW$ , we can calculate the queuing time with:

$$\frac{x * BW * (r_{tt} + q_{t})}{0.99 * BW} = q_{t}$$

As a result we get:

$$q_{t} = \frac{r_{tt} * x}{0.99 - x}$$

For example, limiting the allowed premium bandwidth to 33% of the network capacity ( $x = 0.33$ ) results in a maximum increase of latency of  $\frac{r_{tt}}{2}$ .

The buffer size allocated to the WFQ queues should leave a significant amount of space. This could be used on demand by either the best effort class, or the premium class. With this configuration we minimize the impact on handling traditional best effort traffic under congestion. It is important to note that the usage of the overall queue limit is only available with the standard tail drop behavior of the class of queue, and not with activating WRED on WFQ queues.



## 6 Application Example

We have constructed our DS resource manager to support two classes [8] of premium service: a foreground service, for latency- and jitter-sensitive flows (e.g., multimedia streaming and control), and a background service, for long-lived, high bandwidth but latency-insensitive flows (e.g., bulk data transfer operations). The resource manager changes the bandwidth provided to background reservations dynamically as foreground reservations come and go. When the bandwidth provided to a background reservation changes, callbacks are provided to the application to inform it of the change. This strategy allows bulk data transfers to co-exist with multimedia flows. The amount of bandwidth available for background reservations over a particular time period can then be controlled via policy mechanisms. Our prototype supports multiple foreground reservations but initially only a single background reservation; the extensions required to support multiple background flows are not complex.

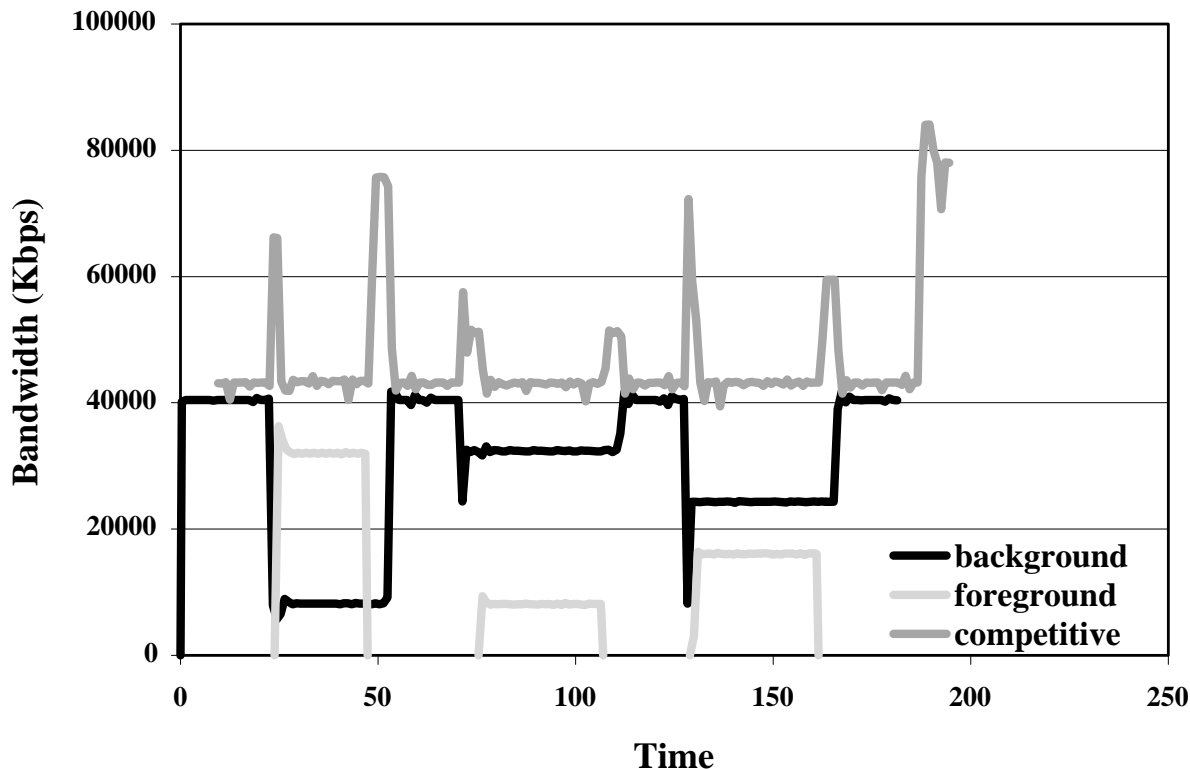


Figure 9: Performance achieved for a mixture of premium and best effort services on GAR-NET. We demonstrate that a bulk-transfer (background) application is able to exploit unused premium traffic without affecting foreground reservations. See text for details.

Figure 9 shows the bandwidth achieved by the foreground, background, and best effort flows during the experiment. We configured GARNET to create a 45 Mb/s premium channel in a 100 Mb/s network. We then created five distinct flows: a bulk data transfer, operating as a “background” flow; a competing 80 Mb/s best-effort UDP flow (a traffic generator submitting 1,000 byte packets every 100  $\mu$ secs); and three independent, short-lived foreground flows with immediate reservations. In this experiment we used a simple data transfer program, which adapted its transmission rate to the rate assigned to the background class.

## 7 Related Work

The general problem of QoS implementation and management is receiving increased attention (e.g., [9]). Some groups have investigated the use of DS mechanisms (e.g., [22, 18]) and its impact on TCP. Especially [18] is presenting a quantitative comparison of different DS router mechanisms. However none of the studies has explicitly addressed flows consuming a goodput of 10 MB/s or higher, using commodity router hardware which provides buffer capabilities of up to several Megabytes.

The pan-European research network, and its task force Testing Advanced Networking Technologies (TF-TANT) are currently evaluating DS mechanism using commodity hardware. TF-TANT addresses similar aspects as discussed in this paper but focuses on flows with lower bandwidth requirements. The results have not yet been published but are available on the web (<http://www.cnaf.infn.it/~ferrari/tfng/ds/del-rep1.doc>).

## 8 Conclusions and Future Work

We have presented a quantitative evaluation of DS implementation for high-performance TCP-flows and demonstrated that end-to-end QoS can be delivered to such flows if DS mechanisms are configured carefully. The basic challenge is dealing with the burstiness introduced by TCP’s sliding window. This must be addressed by appropriate policies at the edge routers; these policies must support bursts correlating to the TCP window size. Burstiness also introduces problems on the interior interfaces, because the available bandwidth might be exceeded by accumulating aggregate bursts. For that reason the implementation of the EF PHB should avoid queueing by overprovisioning the guaranteed amount of bandwidth for the premium class as much as possible.

We have shown that CAR and WFQ can be used to implement a DS architecture. Policing at edge routers is done based on the applied bandwidth and the estimated round-trip time. The EF PHB is implemented by guaranteeing 99% of the available bandwidth to the premium class, which minimizes the maximum latency of premium traffic. We have demonstrated that the impact of this configuration on best effort traffic in the absence of premium traffic is negligible.

In future work we plan to improve the proposed implementation by introducing Distributed Traffic Shaping (DTS), and by analyzing threshold dropping for premium flows exceeding its policy limit at the edge. We will also analyze the impact of priority queueing as an alternative to WFQ.

## Acknowledgments

We gratefully acknowledge assistance given by Robert Olsen who provided additional information about Cisco's DS implementation, by Andy Adamson who wrote the UDP traffic generator, and by Jun Wang who assisted in our measurement effort. Numerous discussions with our colleagues Gary Hoo, Bill Johnston, Carl Kesselman, and Steven Tuecke have helped shape our approach to quality of service. We also thank Cisco Systems for an equipment donation that allowed the creation of the GARNET testbed. This work was supported in part by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Advanced Scientific Computing Research, U.S. Department of Energy, under Contract W-31-109-Eng-38; by the Defense Advanced Research Projects Agency under contract N66001-96-C-8523; by the National Science Foundation; and by the NASA Information Power Grid program.

## References

- [1] M. Allman, V. Paxson, and W. Stevens. TCP Congestion Control. *Internet RFC 2581*, 1997.
- [2] S. Blake, D. Black, M. Carlson, M. Davies, Z. Wang, and W. Weiss. An Architecture for Differentiated Services. *Internet RFC 2475*, 1998.
- [3] R. Braden, D. Clark, and S. Shenker. Integrated services in the internet architecture: an overview. RFC 1633, Internet Engineering Task Force, 1994.
- [4] D. Comer. *Internetworking with TCP/IP*. Prentice-Hall International Editions, 1988.
- [5] W. Feng, D. Kandlur, D. Saha, and K. Shin. Adaptive Packet Marking for Providing Differentiated Services in the Internet. In *Proceedings of the International Conference on Network Protocols, Oct. 1998*. 1998.
- [6] I. Foster and C. Kesselman, editors. *The Grid: Blueprint for a Future Computing Infrastructure*. Morgan Kaufmann Publishers, 1999.
- [7] I. Foster, C. Kesselman, C. Lee, R. Lindell, K. Nahrstedt, and A. Roy. A distributed resource management architecture that supports advance reservations and co-allocation. In *Proceedings of the International Workshop on Quality of Service*, pages 27–36, 1999.
- [8] I. Foster, A. Roy, V. Sander, and L. Winkler. End-to-End Quality of Service for High-End Applications. Technical report, Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, 1999. <http://www-fp.globus.org/documentation/papers.html>.
- [9] Roch Guérin and Henning Schulzrinne. Network quality of service. In [6], pages 479–503.
- [10] J. Heinanen, T. Finland, F. Baker, W. Weiss, and J. Wroclawski. Assured Forwarding PHB Group. *Internet RFC 2597*, 1999.

- [11] V. Jacobson and M. Karels. Congestion Avoidance and Control. In *Proceedings of the SIGCOMM '88*. 1988.
- [12] V. Jacobson, K. Nichols, and K. Poduri. An Expedited Forwarding PHB. *Internet RFC 2598*, 1999.
- [13] T. Lakshman, U. Madhow, and B. Suter. Window-based error recovery and flow control with a slow acknowledgement channel: a study of TCP/IP performance. In *Proceedings of the IEEE INFOCOM '97*. 1997.
- [14] M. Mathis, J. Semke, and J. Mahdavi. The Macroscopic Behavior of the TCP Congestion Avoidance Algorithm. In *Proceedings of ACM SIGCOMM, volume 27, number 3*. 1997.
- [15] M. Mathis, J. Semke, and J. Mahdavi. Automatic TCP Buffer Tuning. In *Proceedings of ACM SIGCOMM, volume 28, number 4*. 1998.
- [16] K. Nichols, S. Blake, F. Baker, and D. Black. Definition of the Differentiated Services Field (DS Field) in the ipv4 and ipv6 headers. *Internet RFC 2474*, 1998.
- [17] R. Nitzan and B. Tierney. Experiences with TCP/IP over an ATM OC12 WAN. LBNL Report, Lawrence Berkeley National Laboratory, April 1999.
- [18] Sambit Sahu, Don Towsley, and Kim Kurose. A Quantitative Study of Differentiated Services for the Internet. UMass CMPSCI Technical Report, University of Massachusetts, Sept 1999.
- [19] W. Stevens. TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms. *Internet RFC 2001*, 1997.
- [20] W. Stevens. *TCP/IP Illustrated, Vol.1 The Protocols*. Addison-Wesley, 1997.
- [21] I. Yeom and N. Reddy. Realizing Throughput Guarantees in a Differentiated Services network. In *Proceedings of the ICMCS, June 99*. 1999.
- [22] Ikjun Yeom and A. L. Narasimha Reddy. Modeling TCP behavior in a Differentiated-Services Network. Technical report, TAMU ECE, 1999.

## 9 Vitae

Volker Sander is a visiting scientist at Argonne National Laboratory's Distributed Systems Laboratory (DSL). He has been actively researching issues in metacomputing since 1995 at Forschungszentrum Juelich GmbH. The focus of his current research is Quality of Service, Remote I/O, and Security.

Dr. Ian Foster is Senior Scientist and Associate Director in the Mathematics and Computer Science Division at Argonne National Laboratory, and also Associate Professor of Computer Science at The University of Chicago. His research interests include parallel and distributed computing, and computational science. He has written 4 books and over 100 technical articles in these and related areas.

Alain Roy is a graduate student from the University of Chicago. He has been actively researching issues in Quality of Service since 1997. He current research is with GARA, a system that allows users to make advanced reservations for not only networks, but for CPUs, disks, and even graphics pipelines.

Linda Winkler is a Senior Network Engineer at Argonne National Laboratory's Mathematics and Computer Science Division. She is currently serving as the MREN Technical Director and is a primary member of the STARTAP engineering team. Her focus since 1995 has been in the area of interconnectivity and interoperability of wide-area research networks in support of advanced scientific and engineering applications.